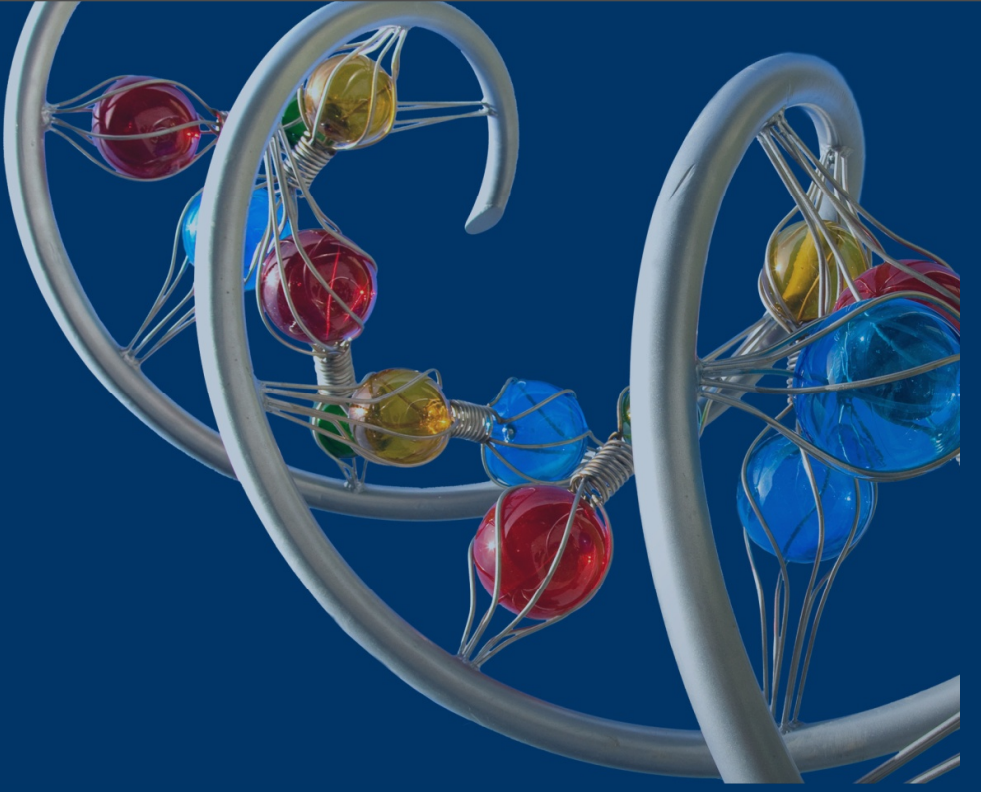


BioPig: Developing Cloud Computing Applications for Next-Generation Sequence Analysis

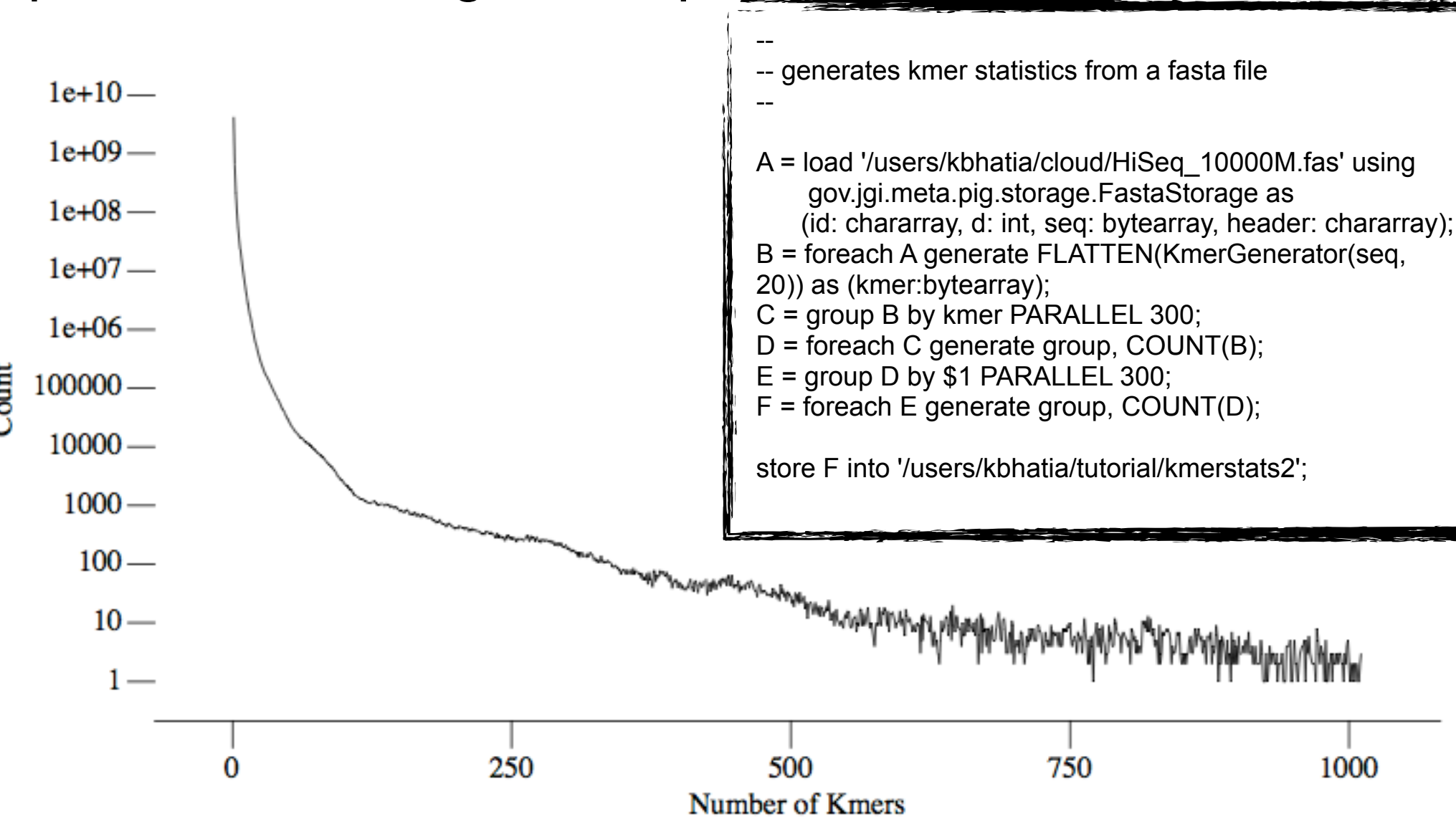


Application

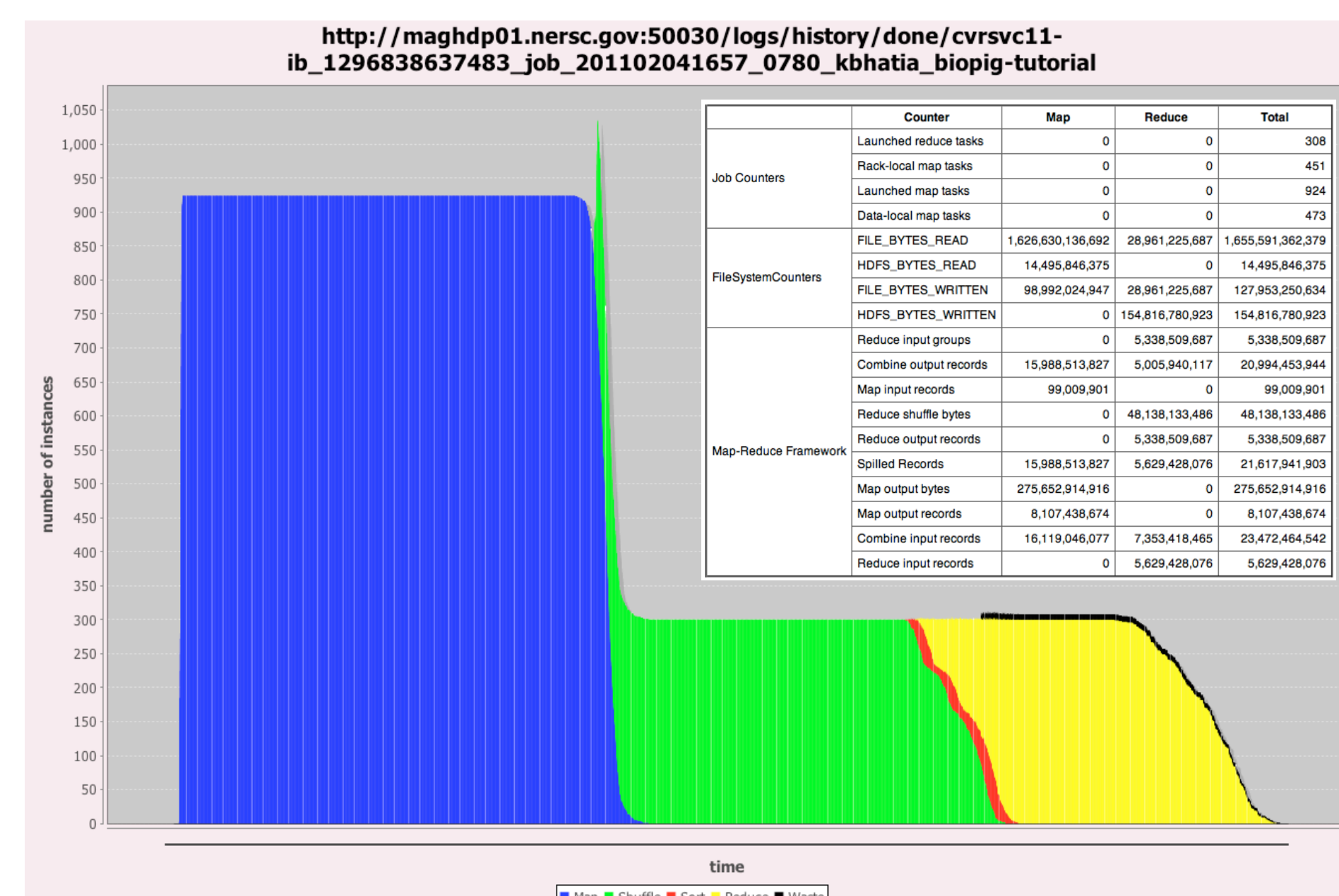
kmer histograms

Kmers (also sometimes called ngrams) represent a sliding window of fixed length (k) across the sequence. Kmers are quite versatile and are the basis for a variety of analytics. A Kmer Histogram is a plot of the frequencies of kmers across a dataset that is an indicator for coverage, the existence of metagenomes or sequencing errors. The below plot shows typical *metagenomic* kmer frequencies for one lane of Illumina HiSeq with $k=20$.

Computationally, Kmer histograms require significant memory to maintain a hash table of counts. However, kmer generation and counting is easily parallelized using Hadoop MapReduce.



Running 1000-way parallelism generates 300GB of intermediate data, 5 Billion kmers, and completes in 15 min.



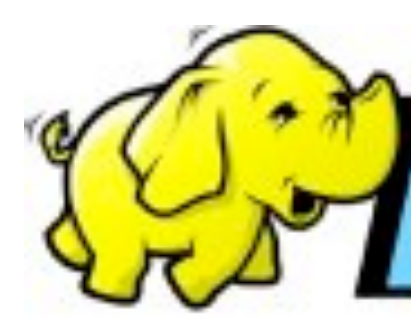
Performance scales linearly with size of sequences, and linearly with cluster size. Double the data, double the nodes.

Amazon Elastic MapReduce provides similar performance with new CC.Large nodes (@\$2/hr each)

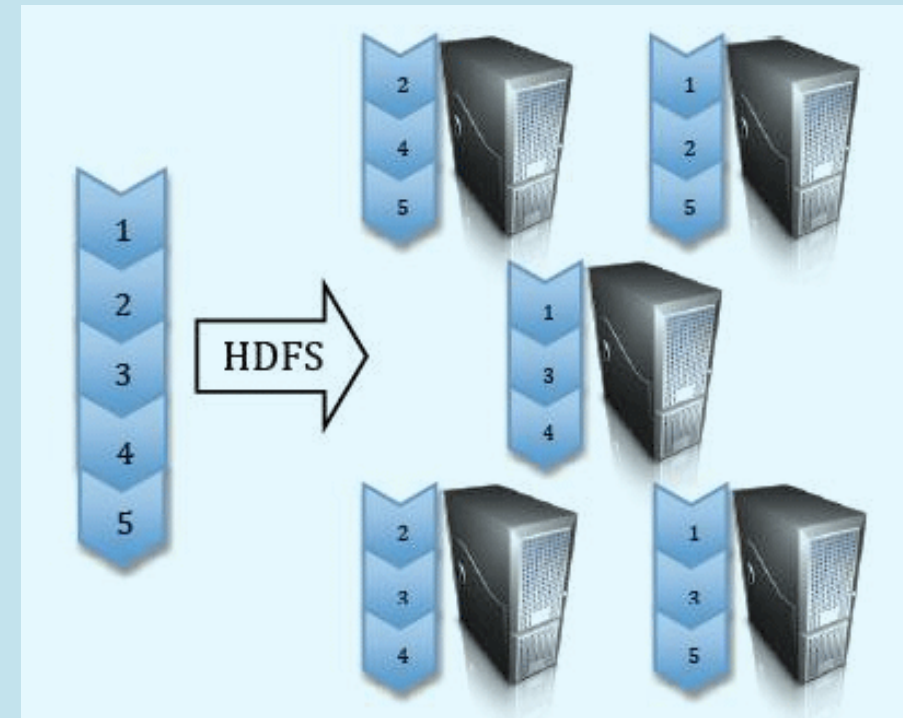
Karan Bhatia, Zhong Wang
Joint Genomics Institute
Lawrence Berkeley National Laboratory

Next Generation sequencing is producing ever larger data sizes with a growth rate outpacing Moore's Law. The data deluge has made many of the current sequence analysis tools obsolete because they do not scale with data. Here we present BioPig, a collection of cloud computing tools to scale data analysis and management. Pig is a flexible data scripting language that uses Apache's Hadoop data structure and map reduce framework to process very large data files in parallel and combine the results. BioPig extends Pig with capability with sequence analysis. We will show the performance of BioPig on a variety of bioinformatics tasks, including screening sequence contaminants, Illumina QA/QC, and gene discovery from metagenome data sets using the Rumen metagenome as an example.

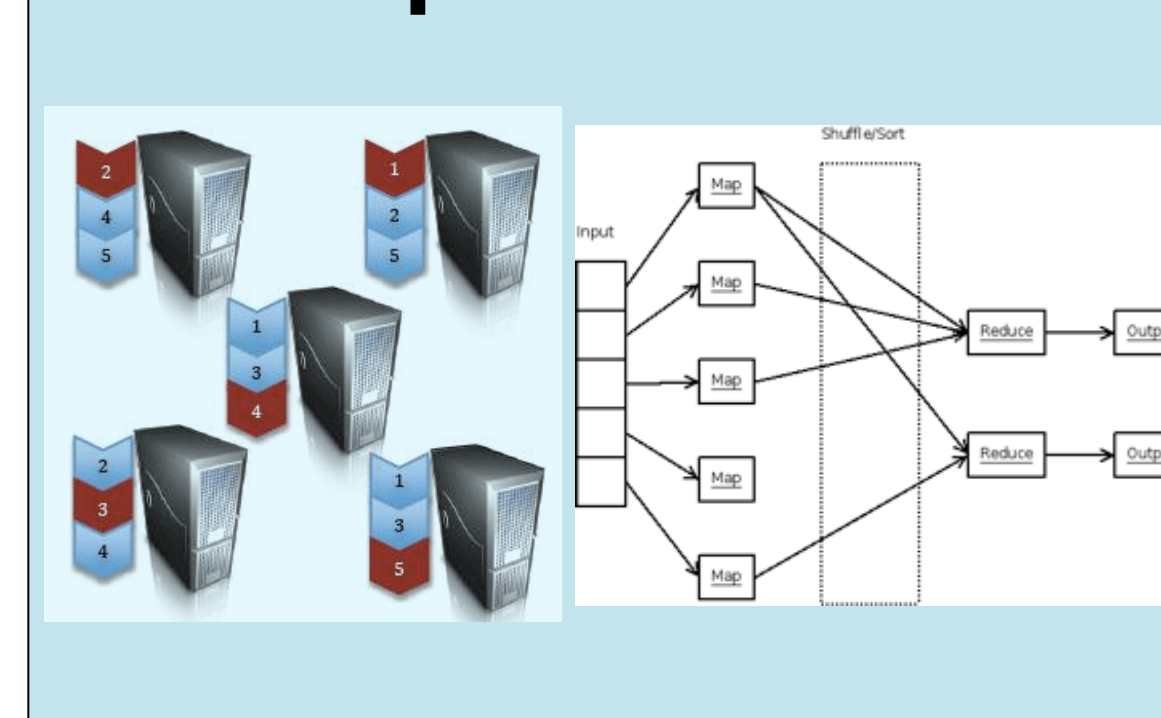
Cloud Technology



HDFS



MapReduce



Pig

Pig Democratizes Large-scale Data Analysis

The Pig version is:
5% of the code
5% of the time
 Within 50% of the execution time.
 Innovation increasingly driven from large-scale data analysis
 Need fast iteration to understand the *right questions*
 More minds contributing = more value from your data

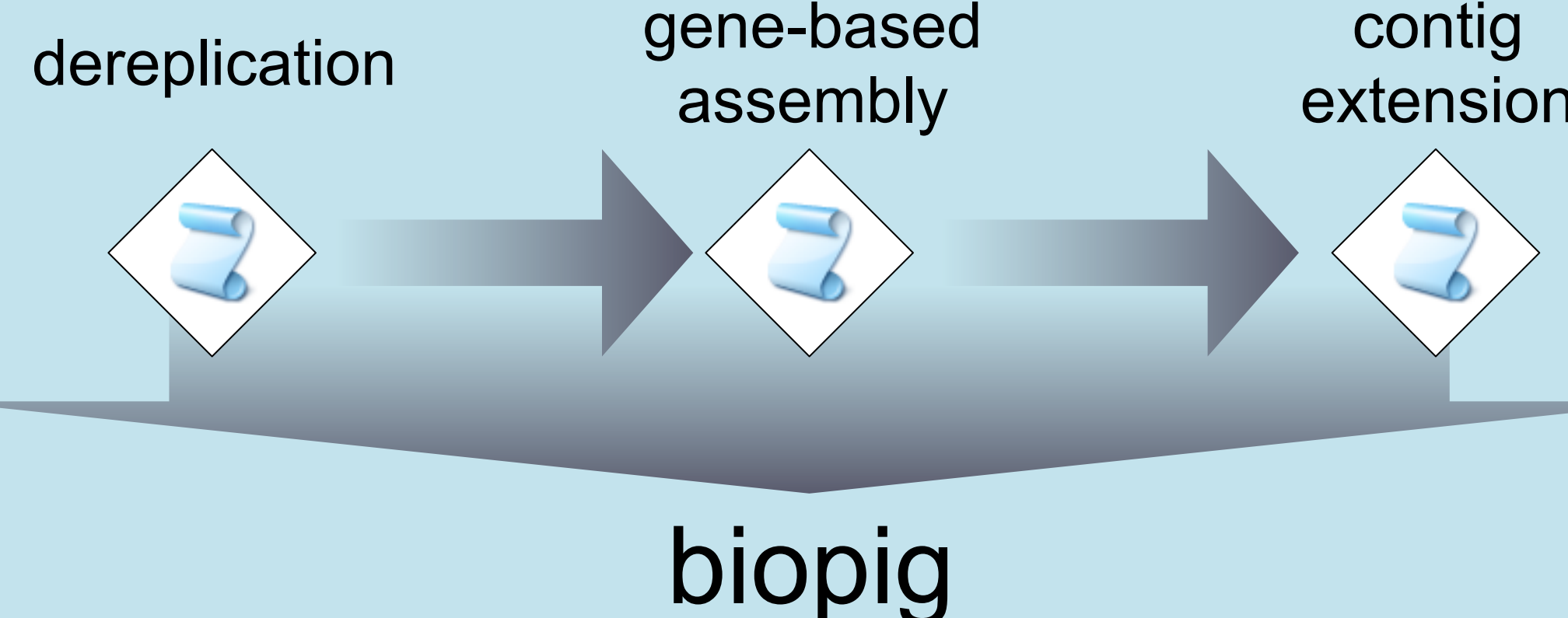
BioPig

features

- Custom data loaders for FASTA and FASTQ data files
- Wrapped many external Bioinformatics applications: Blast, Velvet, Newbler, CAP3
- native support for KMER generation
- supports paired and unpaired sequences
- Additional functions for N50, subsequence sequence neighbors

- leverage all the existing capabilities of Pig including data joins, filters, aggregate, etc.

rumen pipeline



<http://codaset.com/zhongwang/meta>

Application

Screening

Metagenomic sequence data can represent a multitude of individual species found at the sample site, as well as species introduced in the sequencing process. Especially valuable for Single-Cell sequencing, kmer-based pairwise matching to the NT database can identify contaminants, as shown here.



Magellan Cloud Testbed

Base Compute Nodes

- 560 nodes
- 2 quad-core Intel Nehalem 2.67 GHz processors per node
- 8 cores per node (4,480 total cores)
- 24 GB DDR3 1333 MHz memory per node

Expanded Compute Nodes

- 160 nodes - all being used for Magellan research
- 2 quad-core Intel Nehalem 2.67 GHz processors per node
- 8 cores per node (1,280 total cores)
- 48 GB DDR3 1066 MHz memory per node
- 1 TB (local) SATA disk per node

Login/Network Service Nodes

- 18 nodes
- 2 quad-core Intel Nehalem 2.67 GHz processors
- 8 cores per node (144 total cores)
- 48 GB DDR3 1066 MHz memory per node

High Performance Interconnect

- 4X QDR InfiniBand, fibre optic cables
- Local fat-trees with a global 2D mesh

Cooling

- Liquid Cooled



DISCLAIMER

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California.